

3 SYSTEMS AND METHODS FOR SEQUENCE COMPARISON
4

5 CLAIM OF PRIORITY

6 [0001] This application claims priority to U.S.S.N. 60/429,965, entitled “Systems and
7 Methods for Sequence Comparison”, filed on November 29, 2002, the contents of which are
8 herein incorporated by reference in their entirety.

9
10 BACKGROUND

11 (1) Field

12 [0002] The disclosed methods and systems relate generally to pairwise alignment, and more
13 particularly to computing percent identities based on pairwise alignments.

14 (2) Description of Relevant Art

15 [0003] Patent claims are generally directed to a nucleotide and/or polypeptide sequence
16 invention and other sequences having a claimed level (e.g., percentage) of sequence identity to
17 that invention. Freedom to operate and patentability assessments related to such claims can be
18 based on queries of databases of such sequences, where such databases can include, for example,
19 PAT 10, GENESEQ 5, EMBL, GenBank, and DDBJ. Unfortunately, search tools often do not
20 identify sequences that may be within the scope of the patent claims, thus potentially causing an
21 inaccurate evaluation and/or assessment of intellectual property rights.

22 [0004] Patentability and freedom-to-operate assessments can often rely on search and/or
23 query tools that are based on homology. Such tools include, for example, BLAST (Basic Local
24 Alignment Search Tool) and FASTA, which retrieve putative homologues based on a user-
25 provided query sequence. Those of ordinary skill understand that two genes are homologues if
26 they can be understood to have evolved from the same common ancestor gene. BLAST and
27 FASTA can be understood to be directed to identifying a homology relationship between two
28 sequences (e.g., the query sequence, and a database sequence), and accordingly, these tools are

1 based on a set of parameters (e.g., gap opening penalty value, substitution matrix, cut-off scores,
2 gap extension penalty, etc.) and a scoring system that conveys biological information. The
3 configurable parameters can significantly alter the output of a given query, as the queries are
4 based on an internally computed score that can be highly sensitive to these parametric changes.
5 Examples can be provided where the same query may yield a 93% identity for one set of
6 parameters, and a 100% identity with another set.

7 [0005] Particularly, sequence alignments may often be computed with percent identity scores
8 only after the “best hits” are identified by a scoring function. As it’s name implies, BLAST
9 performs and optimizes an alignment with respect to a fraction of the query that gives the highest
10 percentage score. Accordingly, the BLAST alignment can be achieved on a portion of the query,
11 leading to a percent identity only with respect to a fraction of the query. Further, a user generally
12 cannot specify which fraction of the query should be used for the alignment. BLAST can thus
13 report multiple alignment results.

14 [0006] In contrast to BLAST, FASTA provides one alignment result by identifying local high
15 scoring alignments, starting as a BLAST search with exact short word matches and extending
16 potential hits based on a greedy ungapped basis. The result is a single gapped or ungapped
17 alignment result.

18 [0007] Accordingly, regardless of whether BLAST and/or FASTA are employed, both query
19 tools rely on the same homology-based search paradigm and the respective results are
20 additionally based on upon a set of user-provided parameters that can significantly affect the
21 query results. Searching or otherwise querying sequence databases based on patent claims
22 necessitates retrieving database elements that either completely match the query or are related to
23 a specified extent, regardless of homology. The aforementioned tools can thus be considered
24 ineffective for freedom to operate and patentability assessments.

25

26 SUMMARY

27 [0008] The disclosed methods and systems include a method for comparing a first sequence
28 and a second sequence, the method including associating errors with alignments of the first

1 sequence and the second sequence, comparing the alignment errors to identify the alignment
2 having the smallest error, and, based on the alignment having the smallest error, computing a
3 first percent identity relative to the first sequence, and a second percent identity relative to the
4 second sequence. The method can include determining a mismatch number based on mismatches
5 between the first sequence and the second sequence based on the alignment having the smallest
6 error, and/or an alignment number based on matches between the first sequence and the second
7 sequence based on the alignment having the smallest error.

8 [0009] In an embodiment, computing a first percent identity relative to the first sequence can
9 include determining an alignment number based on the matches between the first sequence and
10 the second sequence based on the alignment having the smallest error, and, forming a ratio based
11 on the alignment number and the length of the first sequence. Further, computing a second
12 percent identity relative to the second sequence can include determining an alignment number
13 based on the matches between the first sequence and the second sequence based on the alignment
14 having the smallest error, and, forming a ratio based on the alignment number and the length of
15 the second sequence.

16 [0010] In some embodiments, the methods and systems can include computing a third
17 percent identity relative to the alignment having the smallest error. The third percent identity can
18 be computed by determining an alignment number based on the matches between the first
19 sequence and the second sequence based on the alignment having the smallest error, and, forming
20 a ratio based on the alignment number and the length of the alignment. The matches can be
21 perfect matches and/or positive matches. In an embodiment, a user or another can provide an
22 input as to whether the percent identities can be computed based on perfect and/or positive
23 matches.

24 [0011] In an embodiment, a user or another can provide a percent identity threshold, such
25 that the methods and systems can include determining whether at least one of the first percent
26 identity and the second percent identity is greater than a percent identity threshold.

27 [0012] In some embodiments where the methods and systems can include inserting gaps into
28 either of the first and/or second sequence, the methods and systems can including determining a

1 number based on the gaps in the first sequence based on the alignment having the smallest error,
2 and, a number based on the gaps in the second sequence based on the alignment having the
3 smallest error.

4 [0013] In one embodiment, one or more databases can be provided, where the database(s)
5 can include one or more sequences, and the first and/or second sequences can be retrieved from
6 the database(s). Accordingly, the disclosed methods and systems can allow for a single database
7 of sequences to be compared against itself, and the methods and systems can also be extended to
8 include comparisons of sequences from multiple databases. In one exemplary embodiment, the
9 first sequence(s) can include one or more polypeptide sequence(s) and/or nucleotide sequence(s),
10 and, the second sequence(s) can include one or more polypeptide sequence(s) and nucleotide
11 sequence(s). Accordingly, the methods and systems can allow for serial and/or parallel
12 processing of sequence alignment/comparison using one or more processing threads.

13 [0014] Aligning or otherwise comparing the first sequence and the second sequence can
14 include aligning the first sequence and the second sequence, and computing an error based on the
15 number of mismatches in the alignment. As provided previously herein, the alignment can
16 include one or more insertion events with respect to the first sequence and/or the second
17 sequence. In an embodiment, the alignment can be understood to include computing a string edit
18 distance. In one embodiment, the string edit distance can be associated with the alignment error,
19 and in such an embodiment, the alignment having the smallest error may be associated with the
20 smallest string edit distance for a given first sequence and second sequence.

21 [0015] The disclosed methods and systems can include identifying whether the first sequence
22 is longer than the second sequence, whether the second sequence is longer than the first
23 sequence, and/or whether the first sequence and the second sequence are equivalent and/or equal
24 length. In an embodiment where string edit distance can be determined, the query string can be
25 understood to be the shorter sequence, and the target sequence can be understood to be the longer
26 sequence. Alignment errors can be computed by or otherwise based on determining a string edit
27 distance by computing those alignments where the shorter sequence can be included in (e.g.,
28 overlap) the longer sequence. When the first and second sequences are the same length,

1 alignment errors can be computed based on the first sequence being the shorter sequence, and
2 further, based on the second sequence being the shorter sequence. In one embodiment, the
3 methods and systems can include aligning at least the entirety of the shorter sequence with at
4 least a fragment of the longer sequence. Aligning at least the entirety can include inserting at
5 least one gap into at least one of the shorter sequence and/or the longer sequence. In some
6 embodiments, the alignment can be performed regardless of homology.

7 [0016] Alignment errors, including the alignment having the smallest error, can be compared
8 to an alignment error threshold. A user or another can provide the alignment error threshold, and
9 in some embodiments, an output may be based on the alignment error threshold such that
10 alignments having a number of alignment errors exceeding the alignment error threshold may not
11 be output or otherwise provided, for example.

12 [0017] A user or another can also provide a percent identity threshold, where outputs for the
13 methods and systems can be based on a comparison of the first percent identity and the second
14 percent identity relative to the percent identity threshold. In one example, if either of the first
15 percent identity or the second percent identity exceeds the percent identity threshold, the methods
16 and systems can output (e.g., transmit to display, memory, storage, another application, another
17 device, and/or otherwise provide) data (e.g., first percent identity, second percent identity, third
18 percent identity, scoring matrix(s), scoring matrix(s) metrics (e.g., perfect matches, positive
19 matches, etc.), alignment error data, number of gaps in first sequence, number of gaps in second
20 sequence, alignment identification data, positions of alignments, etc.). Those of ordinary skill
21 can recognize that the methods and systems can include comparing the length of the first
22 sequence with the length of the second sequence, and performing the alignments based on the
23 length comparison and a percent identity threshold. Accordingly, if the length difference(s)
24 between a first and second sequence exceeds a percent identity threshold, alignments may not be
25 performed for the first and second sequence pair.

26 [0018] In an embodiment, the aligning can be performed based on a dynamic programming
27 method for approximate string (e.g., sequence) matching. For example, the methods and systems
28 can include determining locations at which a query string/sequence (e.g., first sequence) of length

1 m matches a sub-string (e.g., subsequence) of a subject string/sequence (e.g., second sequence) of
2 length n, where n is longer than m, and where such locations of alignment can provide less than k
3 errors. K can be an alignment error threshold.

4 [0019] The methods and systems can include one or more interfaces to allow a user or
5 another to identify the first sequence(s) (e.g., database(s)), identify the second sequence(s) (e.g.,
6 database(s)), provide a percent identity threshold, and/or provide an alignment error threshold.
7 Accordingly, the methods and systems can include performing multiple sequence comparisons,
8 and can thus include, iteratively, storing the first percent identity and the second percent identity,
9 retrieving a first sequence and/or a second sequence, and repeating the process of associating
10 errors with the alignments, to provide at least one stored first percent identity and second percent
11 identity, where such stored percent identities can exceed a percent identity threshold, and can be
12 associated with alignments having a number of alignment errors less than or equal to an
13 alignment error threshold. The stored percent identities can be associated with the first and
14 second sequences to which they apply, and such storage can be sorted based on, for example,
15 percent identity.

16 [0020] Other objects and advantages will become apparent hereinafter in view of the
17 specification and drawings.

18

BRIEF DESCRIPTION OF THE DRAWINGS

19 [0021] Figures 1A and 1B depict illustrative embodiments of the disclosed systems and
20 methods;
21 Figure 2 is an exemplary block diagram providing some components of an illustrative
22 system and method;
23 Figures 3A-3C show local alignment, global alignment, and best-fit alignments; and,
24 Figure 4 is an example scoring matrix.

25

DESCRIPTION

1 [0022] To provide an overall understanding, certain illustrative embodiments will now be
2 described; however, it will be understood by one of ordinary skill in the art that the systems and
3 methods described herein can be adapted and modified to provide systems and methods for other
4 suitable applications and that other additions and modifications can be made without departing
5 from the scope of the systems and methods described herein.

6 [0023] Unless otherwise specified, the illustrated embodiments can be understood as
7 providing exemplary features of varying detail of certain embodiments, and therefore, unless
8 otherwise specified, features, components, modules, and/or aspects of the illustrations can be
9 otherwise combined, separated, interchanged, and/or rearranged without departing from the
10 disclosed systems or methods. Additionally, the shapes and sizes of components are also
11 exemplary and unless otherwise specified, can be altered without affecting the scope of the
12 disclosed and exemplary systems or methods of the present disclosure.

13 [0024] The disclosed methods and systems include methods and systems to compare two
14 sequences that include a first sequence and a second sequence. In one embodiment, the
15 sequences can be polypeptide and/or nucleotide sequences, although the methods and systems are
16 not so limited, and other sequences of ASCII characters can be used, where such sequences can
17 apply to other applications and/or embodiments. Generally, references to the word "sequence"
18 herein can be understood to be an ASCII string. The methods and systems can be used to
19 compare the lengths of the first and second sequences to determine a shorter sequence and a
20 longer sequence, and to determine a best fit of the shorter sequence in the longer sequence.
21 Based at least on the best fit, a first percent identity can be computed relative to the first
22 sequence, and a second percent identity can be computed relative to the second sequence. The
23 first and second percent identities can be provided as an output to a display, a database, a
24 memory, a computer program, another processor-controlled device, and/or another output device.
25 In one embodiment, where the first and second sequences can be based on polypeptide and/or
26 nucleotide sequences, the methods and systems can be employed to determine whether the first
27 and second sequences may be within the scope of a patent claim that may be associated with one
28 or both of the first sequence and the second sequence.

1 [0025] In one embodiment, a first sequence can be a sequence provided in and/or otherwise
2 proposed for inclusion in a patent application, while the second sequence(s) can be one or more
3 prior art sequence, and thus, the methods and systems may be applied to the first and second
4 sequence(s) to determine whether the first sequence is novel when compared to the second
5 sequence(s), and/or whether an entity is free to use such first sequence (e.g., a freedom to operate
6 analysis as is known in the art). As provided herein, the first and/or second sequence can be
7 understood more generally to be a string, such as a biological sequence and/or a synthetic
8 sequence, with such examples provided for illustration and not limitation.

9 [0026] Figure 1A shows one illustrative embodiment for implementing the disclosed
10 methods and systems. As Figure 1 indicates, a user-device **10**, which can be a processor-
11 controlled device as provided herein, can be provided with a user-interface **12** that can allow a
12 user or another to input information and/or data that can be used by the disclosed methods and
13 systems. As Figure 1A shows, such information can include an identity of a first sequence and
14 an identity of a second sequence. Those of ordinary skill will understand that in one
15 embodiment, the first sequence can be, for example, a “query” sequence, while the second
16 sequence can be and/or include one or more databases which can include or otherwise contain
17 one or more “target” sequences to which the first/query sequence can be compared. In other
18 embodiments, the first sequence and the second sequence can be individual sequences, and the
19 first sequence can be associated with the target, while the second sequence can represent the
20 query. In one embodiment, the first sequence can be a database and/or otherwise associated
21 therewith, and the second sequence can be a database and/or otherwise associated therewith, and
22 in some instances, such databases can be the same to indicate a request for performing the
23 methods and systems against a single database. In other embodiments, the databases can be
24 different. Accordingly, those of ordinary skill will recognize that references herein to first
25 sequence and second sequence can include one or more identifiers for one or more query
26 sequences, and/or one or more identifiers for one or more target sequences, such as one or more
27 databases of target sequences. Such identifiers can be provided via a user input or other
28 designation means.

1 [0027] Accordingly, the designations and/or references to first and second sequences herein
2 can be understood to be arbitrary, and for the discussion herein, it can be understood that the first
3 sequence is the query sequence (or a reference and/or identifier related thereto, e.g., one or more
4 databases), while the second sequence is the target sequence (or a reference and/or identifier
5 related thereto, e.g., one or more databases).

6 [0028] Referring again to Figure 1A, in some embodiments, a user or another can provide,
7 via a user interface or other means (e.g., specified by a system manager, hard-coded) a percent
8 identity threshold which can indicate a percent identity of comparison between the first and
9 second sequences, such that the computed percent identities relative to the first sequence and the
10 second sequence can be compared to the percent identity threshold. In one embodiment, for
11 example, computed percent identities may be stored when one or both of such computed percent
12 identities equal or exceed the threshold. Other comparisons to the threshold can be made.

13 [0029] Figure 1A also indicates that the user or another can enter, via a user interface or
14 another means, an alignment error threshold that can be used, for example, for comparison
15 against a number of alignment errors between two sequences. In an embodiment, if a number of
16 errors between the two sequences (first sequence and second sequence) equal or exceed the
17 alignment error threshold, such alignment may not be deemed to be desirable by the user, may
18 not be output by the system, and/or may not otherwise be included in computations, where such
19 examples are provided for illustration and not limitation.

20 [0030] Those of ordinary skill will recognize that the use of a percent identity threshold and
21 an alignment error threshold can be optional, and as provided herein previously, although in
22 some embodiments, such thresholds can be variable and/or user-specified, in other embodiments,
23 the thresholds may be fixed by, for example, a system administrator, user, or another.

24 [0031] As indicated by Figure 1A, the user entered information can be provided to one or
25 more servers 16, where such servers 16 can be understood to be associated with one or more
26 processor controlled devices as provided herein. Such servers 16 can include instructions for
27 accepting the user-provided information and for accessing processor-executable instructions as
28 provided herein for providing and/or otherwise performing sequence comparisons/alignments.

1 The servers **16** can have access to one or more databases **18** which can include databases of
2 sequences such as polypeptide and/or nucleotide sequences, although other sequences can be
3 included, and such sequences can include alphanumeric sequences, binary sequences, and/or
4 other strings. In an embodiment according to Figure 1A, the user can request a comparison by
5 providing the aforementioned user-specified information at a user device **10**, where such
6 information can be transmitted to a server(s) **16** via a wired or wireless connection using one or
7 more intranets and/or the internet, where the servers **16** can thereafter process the request by
8 accessing the databases **18**. Such database accessing can include querying the databases **18** based
9 on the user information. Upon completing the requested sequence comparisons, the servers **16**
10 can provide the user-device **10** with outputs and/or results that can be provided to a memory, the
11 device display **14**, or other location.

12 [0032] Those of ordinary skill in the art will recognize that the Figure 1A illustrative system
13 can be understood to be representative of a client-server paradigm, where the instructions on the
14 user device **10** for obtaining user information and requesting a comparison can be a client, and
15 the servers **16** can be a server in the client-server paradigm. Accordingly, it can be understood
16 that the illustrated user device **10** instructions and instructions on the servers **16** can be included
17 in a single device such as provided by Figure 1B, where such embodiment may also be
18 considered within the client-server paradigm. As Figure 1B indicates, the user device **10** can
19 access, via wired or wireless communications and using one or more intranets and/or the internet,
20 the databases **18** for querying and/or retrieving sequences. Additionally, the Figure 1B
21 embodiment can represent an embodiment that may not include a client-server paradigm.

22 [0033] Referring to Figure 2, there is one illustrative block diagram for one embodiment of
23 the disclosed methods and systems. As Figure 2 indicates, using a user device **10** or another
24 means, a user or another can provide an identifier to identify at least one first sequence, at least
25 one second sequence, and optionally, a percent identity threshold, and an alignment error
26 threshold **20**. For the purposes of discussion with respect to the illustrative embodiments,
27 reference can be made to a single first sequence and a single second sequence, although it can be
28 understood as provided previously herein, that the methods and systems can be applied to one or

1 more first sequences, one or more second sequences, where such sequences can be in a single
2 and/or multiple databases, and thus such discussion is merely for convenience and can be
3 understood to encompass or otherwise embody multiple first and/or second sequences.

4 [0034] Referring again to Figure 2, the first and second sequences can be retrieved or
5 otherwise provided, obtained, and/or identified 22, whereupon the lengths of the sequences can
6 be determined and compared 24. In one embodiment, such as the embodiment according to
7 Figure 2, for example, if the comparison of the lengths is less than the percent identity threshold
8 26, no further processing may be performed, and the disclosed method and system can retrieve
9 another first sequence and/or second sequence 22 as provided by the embodiment. For example,
10 if the percent identity threshold is ninety-percent, the exemplary embodiment requests alignments
11 equal to or exceeding such percent identity threshold, and the first sequence is length ten, while
12 the second sequence is length one-thousand, one of ordinary skill in the art will recognize that
13 such length comparison can indicate that such percent threshold cannot be satisfied using an
14 alignment, and thus, determining an alignment may not be performed. Such a decision 26 can be
15 optional, as may other decisions in other illustrative embodiments.

16 [0035] According to the Figure 2 embodiment, when the relative and/or comparative
17 sequence lengths equals and/or exceeds the percent identity threshold 26, one of the first or
18 second sequences can be identified as the shorter (“query”) sequence, while the other such
19 sequence can be identified as the longer (“target”) sequence 28. In embodiments when the first
20 sequence and the second sequence may be the same length, the methods and systems disclosed
21 herein can be applied to such sequences with the first sequence identified as the shorter sequence
22 and the second sequence identified as the longer sequence, and again, with the first sequence
23 identified as the longer sequence and the second sequence identified as the shorter sequence. As
24 provided previously herein, a best-fit module can be applied to the sequences 30 to determine
25 and/or identify at least one fragment of the longer sequence which can provide a best-fit 33 for at
26 least the entirety of the shorter sequence. The best-fit can be understood to include one or more
27 alignments. The best fit can be determined based on computing a string edit distance, which can

1 be computed in one or more manners, and such computation can additionally and/or optionally
2 include other dynamic programming methods and systems.

3 [0036] Accordingly, in one embodiment, the disclosed methods and systems can be
4 understood to include computing, determining, or otherwise providing an edit distance between
5 two strings and/or sequences (“string edit distance”), where the string edit distance can be
6 understood to be a minimum number of character inserts (“insertion event”) and/or changes to
7 convert a first string/sequence, to a second string/sequence, where the second sequence can be
8 greater than or equal in length to the first sequence. The insertion events can occur in either the
9 first string and/or the second string. For example, given a first sequence of length m (“pattern”
10 sequence), and a second sequence of length n, informally, the string/sequence edit distance can
11 include computing the smallest edit distance between the first sequence and sub-strings of the
12 second sequence. A sub-string matching method and/or system can thus be understood to
13 identify a best-fit position of a given sub-string (e.g., shorter of the first sequence and the second
14 sequence) within another longer, string (e.g., longer of the first sequence and the second
15 sequence). In one embodiment, a result can include the beginning position within the target
16 (longer) string/sequence where the best match (e.g., minimum number of errors between the
17 sequences) is found.

18 [0037] Those of ordinary skill will recognize that the “best-fit” methods and systems can be
19 distinguished from local alignment methods and systems that may, for example, align only a
20 portion (e.g., sub-sequence) of the shorter (“query”) sequence/string, and similarly, can be
21 distinguished from global alignment methods and systems that may, for example, attempt to align
22 the entirety of the shorter (“query”) sequence/string to the entirety of the longer (“target”)
23 sequence/string. The disclosed methods and systems, as provided herein, attempt to align at least
24 the entirety of the shorter sequence to at least a fragment of the longer sequence. Those of
25 ordinary skill will understand that aligning at least the entirety of the shorter sequence may
26 include an alignment such that a portion of the shorter sequence aligns with at least a fragment of
27 the longer sequence, where such at least a fragment of the longer sequence may include one or

1 both of the ends of the longer sequence. An example of the local alignment, global alignment,
2 and best-fit alignment is shown in Figures 3A-C, respectively.

3 [0038] Referring again to Figure 2, for at least the alignment having a smallest string edit
4 distance compared to the other alignments and associated string edit distances (“alignment
5 having the smallest error”), the methods and systems can at least compute a first percent identity
6 relative to the first sequence and a second percent identity relative to the second sequence 32.
7 Such computations can include determining an alignment number, where the alignment number
8 can be based on the matches between the first sequence and the second sequence, and such
9 matches can be based on the alignment having the smallest error. Those with ordinary skill will
10 understand that such alignment can include “gaps” in the first sequence and/or the second
11 sequence, and thus the length of the alignment can be longer than, for example, the shorter
12 sequence, as such alignment can include gaps in the shorter sequence.

13 [0039] In one embodiment, the disclosed methods and systems can compute an alignment
14 number based on the number of perfect matches in the alignment, and/or the number of positive
15 matches in the alignment. Those of ordinary skill will understand a positive alignment to include
16 an acceptable substitution, such as when the first and second sequences include amino acids, and
17 one or more amino acids can mutate into another amino acid to allow a positive, rather than a
18 perfect, match. In an embodiment, a user or another can provide an input as to whether the
19 percent identities can be computed based on perfect and/or positive matches.

20 [0040] The disclosed methods and systems can thus include computing percent identities
21 based on the number of perfect matches relative to the length of the first sequence, the length of
22 the second sequence, and/or the length of the alignment having the smallest error. The methods
23 and systems can also include computing percent identities based on the number of positive
24 matches relative to the length of the first sequence, the length of the second sequence, and/or the
25 length of the alignment having the smallest error. The methods and systems can include
26 computing percent identities based on the number of perfect and positive matches relative to the
27 length of the first sequence, the length of the second sequence, and/or the length of the alignment
28 having the smallest error. Those of ordinary skill will recognize that a percent identity can

1 include forming a ratio based on one or more of the aforementioned alignment numbers (e.g.,
2 negative, perfect, positive matches), and a length of one or more of the different sequences.

3 [0041] Accordingly, although not shown in the Figure 2 embodiment, as provided herein, a
4 third percent identity can be computed by determining an alignment number based on the
5 matches between the first sequence and the second sequence based on the alignment having the
6 smallest error, and, forming a ratio based on the alignment number and the length of the
7 alignment.

8 [0042] In an embodiment not shown in Figure 2, a decision may be additionally and/or
9 optionally based on the aforementioned alignment error threshold. Accordingly, the methods and
10 systems can include, for example, determining a number (e.g., an error) based on the number of
11 mismatches based on the alignment having the smallest error, and if such number of mismatches
12 exceeds the alignment error threshold, the methods and systems may return to retrieving another
13 first and/or second sequence 22.

14 [0043] Referring again to Figure 2, based on the percent identities computed relative to the
15 first and second sequences, and whether at least one of such percent identities equals or exceeds
16 the percent identity threshold, and/or whether an alignment number is less than or equal to an
17 alignment error threshold, a scoring matrix can be computed 34 for the at least one best-fit
18 alignment. Those of ordinary skill will recognize that scoring matrices can include the first
19 sequence, the second sequence, and indicators identifying, for each element in the respective
20 sequence, whether there is a perfect match (“|”), a positive match (“+”), or a negative match (“ ”).
21 Such indicators are exemplary. The scoring matrix can be based on the alignment have the
22 smallest error. One example of a scoring matrix is shown in Figure 4.

23 [0044] The disclosed methods and systems can include computing a number based on the
24 gaps in the first sequence, and a number based on the gaps in the second sequence, for an
25 alignment having a smallest error. Further, those of ordinary skill will understand that a “best-
26 fit” alignment may include multiple alignments for a given first and second sequence, where such
27 multiple “best-fit” alignments thus provide different alignment results.

1 [0045] In the Figure 2 embodiment, the scoring matrix data provided herein (e.g., scoring
2 matrix, number of “perfect”, “positive”, and “negative” matches (e.g., “mismatch”), number of
3 gaps in first and/or second sequence, etc., and the aforementioned percent identities can be
4 optionally stored 36 and the methods and systems can retrieve another first and/or second
5 sequence 22. The alignment and/or comparison data that is stored can be optionally sorted 38
6 and/or otherwise organized and output 40 to a user or another, to provide one or more data files
7 that can include graphics and/or other data, where such output data can be transmitted 40 via
8 wired and/or wireless networks to, for example, the user device 10, a display, a memory, another
9 processor-controlled device, and/or another specified location. The output data can be based on
10 one or more of the aforementioned quantities (e.g., percent identities, numbers, scoring matrix
11 metrics, etc.) and/or data items (e.g., scoring matrix), although the illustrated embodiment of
12 Figure 2 includes a subset of such output data. In one embodiment, a user can select output
13 formats and/or locations, and can input data, for example, to indicate whether a scoring matrix
14 data and/or metrics can be computed and/or output.

15 [0046] The methods and systems described herein are not limited to a particular hardware or
16 software configuration, and may find applicability in many computing or processing
17 environments. The methods and systems can be implemented in hardware or software, or a
18 combination of hardware and software. The methods and systems can be implemented in one or
19 more computer programs, where a computer program can be understood to include one or more
20 processor executable instructions. The computer program(s) can execute on one or more
21 programmable processors, and can be stored on one or more storage medium readable by the
22 processor (including volatile and non-volatile memory and/or storage elements), one or more
23 input devices, and/or one or more output devices. The processor thus can access one or more
24 input devices to obtain input data, and can access one or more output devices to communicate
25 output data. The input and/or output devices can include one or more of the following: Random
26 Access Memory (RAM), Redundant Array of Independent Disks (RAID), floppy drive, CD,
27 DVD, magnetic disk, internal hard drive, external hard drive, memory stick, or other storage

1 device capable of being accessed by a processor as provided herein, where such aforementioned
2 examples are not exhaustive, and are for illustration and not limitation.

3 [0047] The computer program(s) can be implemented using one or more high level
4 procedural or object-oriented programming languages to communicate with a computer system;
5 however, the program(s) can be implemented in assembly or machine language, if desired. The
6 language can be compiled or interpreted.

7 [0048] As provided herein, the processor(s) can thus be embedded in one or more devices
8 that can be operated independently or together in a networked environment, where the network
9 can include, for example, a Local Area Network (LAN), wide area network (WAN), and/or can
10 include an intranet and/or the internet and/or another network. The network(s) can be wired or
11 wireless or a combination thereof and can use one or more communications protocols to facilitate
12 communications between the different processors. The processors can be configured for
13 distributed processing and can utilize, in some embodiments, a client-server model as needed.
14 Accordingly, the methods and systems can utilize multiple processors and/or processor devices,
15 and the processor instructions can be divided amongst such single or multiple processor/devices.

16 [0049] The device(s) or computer systems that integrate with the processor(s) can include,
17 for example, a personal computer(s), workstation (e.g., Sun, HP), personal digital assistant
18 (PDA), handheld device such as cellular telephone, laptop, handheld, or another device capable
19 of being integrated with a processor(s) that can operate as provided herein. Accordingly, the
20 devices provided herein are not exhaustive and are provided for illustration and not limitation.

21 [0050] References to “a microprocessor” and “a processor”, or “the microprocessor” and “the
22 processor,” can be understood to include one or more microprocessors that can communicate in a
23 stand-alone and/or a distributed environment(s), and can thus be configured to communicate
24 via wired or wireless communications with other processors, where such one or more processor
25 can be configured to operate on one or more processor-controlled devices that can be similar or
26 different devices. Use of such “microprocessor” or “processor” terminology can thus also be
27 understood to include a central processing unit, an arithmetic logic unit, an application-specific

1 integrated circuit (IC), and/or a task engine, with such examples provided for illustration and not
2 limitation.

3 [0051] Furthermore, references to memory, unless otherwise specified, can include one or
4 more processor-readable and accessible memory elements and/or components that can be internal
5 to the processor-controlled device, external to the processor-controlled device, and/or can be
6 accessed via a wired or wireless network using a variety of communications protocols, and unless
7 otherwise specified, can be arranged to include a combination of external and internal memory
8 devices, where such memory can be contiguous and/or partitioned based on the application.
9 Accordingly, references to a database can be understood to include one or more memory
10 associations, where such references can include commercially available database products (e.g.,
11 SQL, Informix, Oracle) and also proprietary databases, and may also include other structures for
12 associating memory such as links, queues, graphs, trees, with such structures provided for
13 illustration and not limitation.

14 [0052] References to a network, unless provided otherwise, can include one or more intranets
15 and/or the internet. References herein to microprocessor instructions or microprocessor-
16 executable instructions, in accordance with the above, can be understood to include
17 programmable hardware.

18 [0053] Unless otherwise stated, use of the word “substantially” can be construed to include a
19 precise relationship, condition, arrangement, orientation, and/or other characteristic, and
20 deviations thereof as understood by one of ordinary skill in the art, to the extent that such
21 deviations do not materially affect the disclosed methods and systems.

22 [0054] Throughout the entirety of the present disclosure, use of the articles “a” or “an” to
23 modify a noun can be understood to be used for convenience and to include one, or more than
24 one of the modified noun, unless otherwise specifically stated.

25 [0055] Elements, components, modules, and/or parts thereof that are described and/or
26 otherwise portrayed through the figures to communicate with, be associated with, and/or be
27 based on, something else, can be understood to so communicate, be associated with, and or be
28 based on in a direct and/or indirect manner, unless otherwise stipulated herein.

1 **[0056]** Although the methods and systems have been described relative to a specific
2 embodiment thereof, they are not so limited. Obviously many modifications and variations may
3 become apparent in light of the above teachings. Many additional changes in the details,
4 materials, and arrangement of parts, herein described and illustrated, can be made by those
5 skilled in the art. Accordingly, it will be understood that the following claims are not to be
6 limited to the embodiments disclosed herein, can include practices otherwise than specifically
7 described, and are to be interpreted as broadly as allowed under the law.